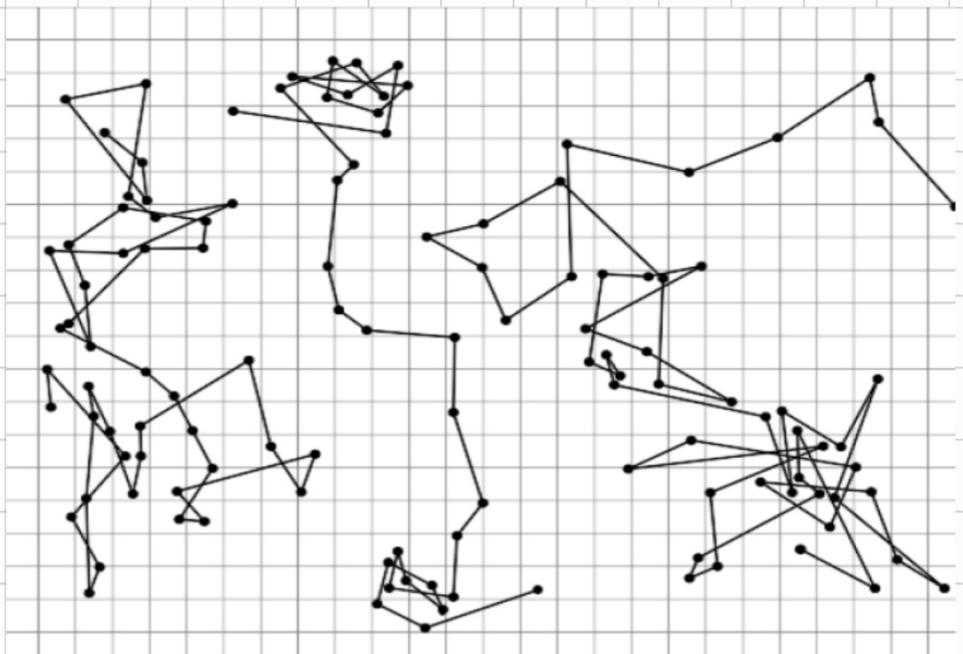


- * Examples :
 - 1) Does my data come from a prescribed distribution F_X ? "Goodness of fit"
 - 2) Rolling a 6-sided die n times and observe 1, 3, 1, 6, 4, ... Is it a fair die?
 - 3) Einstein's Theory of Brownian Motion



* Einstein's Theory of Brownian Motion

P_t : position of the particular at time t .

$P_{t+\Delta t}$: position of the particular at time $t + \Delta t$

$$P_{t+\Delta t} \sim N \left(P_t, \begin{pmatrix} \alpha^2 & 0 \\ 0 & \alpha^2 \end{pmatrix} \right) \text{ where } \alpha^2 = \frac{RT(\Delta t)}{3\pi\eta r N_A}$$

We want to test the goodness of fit.

Experiment : Measured the position of a particle every 30s

1. write T , η , r , N_A and α^2 in terms of σ^2

To verify Einstein's theory and compare it with NA.

$$\text{He figured out } \alpha^2 = 2.23 \times 10^{-7} \text{ cm}^2$$

Question : Does Perrin's data fit with Einstein's model ? ToH - Question.

* Def : Hypothesis Test

It is a binary question about the data distribution. Our goal is to either accept a null hypothesis (H_0 : specifies something about the distribution) or to reject null hypothesis H_0 in favour of an alternative hypothesis H_1 .

If H_0 (similarly H_1) completely specifies the probability distribution for the data then, it is called a simple hypothesis. Otherwise, it is called composite hypothesis.

* Rolled a die n times and observe the outcome.
 X_i denotes the number of time we obtain i , $i \in \{1, \dots, 6\}$

$$H_0 : X_i \sim \text{Multinomial}(n, (\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}))$$

$$H_1 : X_i \sim \text{Multinomial}(n, (\frac{1}{9}, \frac{1}{9}, \frac{2}{9}, \frac{2}{9}, \frac{1}{9}, \frac{2}{9}))$$

Both are simple hypothesis.

Composite alternative hypothesis :

$$H_1 : X_i \sim \text{Multinomial}(n, (p_1, \dots, p_6)) ; p_i \neq \frac{1}{6} \quad \forall i$$

* Another problem is creating a Test Statistic.

$T := T(X_1, \dots, X_n)$ is any statistic; extreme values (large / small) of T provide evidence against H_0 .

$$T = \left(\frac{X_1}{n} - \frac{1}{6} \right)^2 + \dots + \left(\frac{X_n}{n} - \frac{1}{6} \right)^2.$$

Large of T provide evidence against the null hypothesis of a fair die.

Going Back to Brownian Motion:

Let $(X_1, Y_1), (X_2, Y_2), \dots$ be the displacement vector $P_{z_0} - P_0, P_{z_0} - P_{z_0}, \dots$ where $P_t \in \mathbb{R}^2$ is the position of the particle at time t in Brown's experiment.

Simple : $H_0: (X_i, Y_i) \stackrel{iid}{\sim} N(0, 2.23 \times 10^{-7} I)$

Composite : $H_1: (X_i, Y_i) \stackrel{iid}{\sim} N(0, \Sigma I)$ for some $\Sigma \neq 0$

$$\text{Test Statistic : } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$V = \frac{1}{n} \sum_{i=1}^n (X_i^2 + Y_i^2)$$

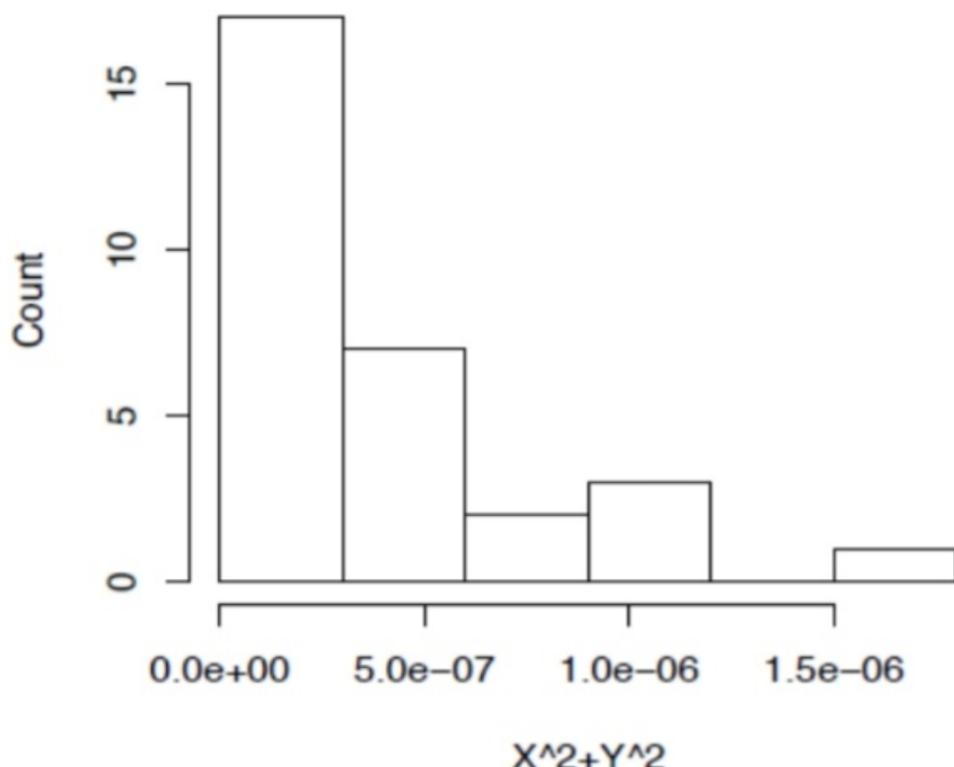
Large / smaller values than 0 for \bar{X} and \bar{Y} provide evidence against H_0 . Values larger / smaller than 2.23×10^{-7} provide evidence against H_0 .

We can rewrite : $V = \sum_{i=1}^n R_i$, $R_i = X_i^2 + Y_i^2$ where $R_i \sim 2.23 \times 10^{-7} \chi_2^2$ (as we have two variables)

* How to derive test statistic from hypothesis ;

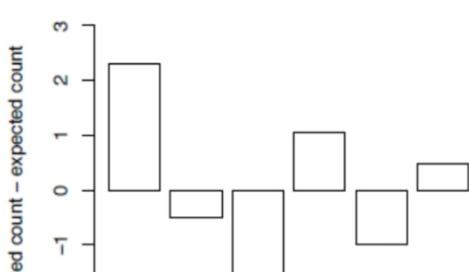
let $R_i = X_i^2 + Y_i^2$. We are interested in testing
 R_1, \dots, R_n are distributed as $\alpha \cdot 2.3 \times 10^{-7} \chi^2_2$ distribution
(under H_0).

Histogram of X^2+Y^2



Deviations from $\alpha \cdot \chi^2$ are better visualized by a hanging histogram which plots $O_i - E_i$ where O_i is the observed count for bin i and E_i is the expected count for bin i under the H_0 distribution.

Hanging histogram of X^2+Y^2



Observ

6
5
4
3
2
1
0

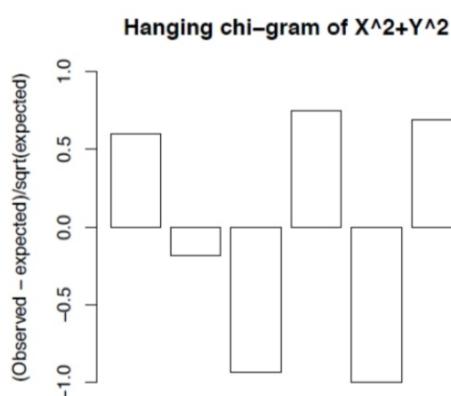
$$\text{A test statistic can be } T = \sum_{i=1}^6 (O_i - E_i)^2.$$

Let p_i be the probability that the hypothesized χ^2 dist. assigns to bin i .

If H_0 is true then $O_i \sim \text{Bin}(n, p_i)$ so $E_i = \mathbb{E}(O_i) = np_i$ and $V(O_i) = np_i(1-p_i) = \mathbb{E}(O_i - E_i)^2$.

The variation in O_i is smaller and it scales approximately linearly with p_i if p_i is close to 0.

We want to stabilize the variance : $\frac{O_i - E_i}{\sqrt{E_i}} = \frac{O_i - E_i}{\sqrt{np_i}}$



The test statistic $T = \sum_{i=1}^6 \frac{(O_i - E_i)^2}{E_i}$ is called **Pearson's chi-squared statistic for goodness of fit.**³

* Test statistic for Q-Q plots :

A Q-Q plot (or probability plot) compares the sorted

values of R_i with $\frac{1}{n+1}, \frac{2}{n+1}, \dots, \frac{n}{n+1}$ quantiles of the hypothesised $\alpha^2 \chi^2$ distribution.

Values close to the line $y = n$ indicate a good fit. F is the cdf of $\alpha^2 \chi^2$ distribution $|R_{(i)} - F^{-1}\left(\frac{i}{n+1}\right)|$.

One way is to take the maximum vertical deviations from the $y = n$ line.

Let $R_{(1)}, \dots, R_{(n)}$ be the ordered statistic of R_1, \dots, R_n . Then, we consider $T = \max_{1 \leq i \leq n} |R_{(i)} - F^{-1}\left(\frac{i}{n+1}\right)|$ where $F^{-1}(t)$ is the t^{th} quantile.

For the values of R where the distribution has high density, the quantiles are closer together so we expect a smaller vertical deviation. This explain why we see more vertical deviations in the upper right of the Q-Q plot.

Another test statistic: $T = \max_{1 \leq i \leq n} |F(R_i) - \frac{i}{n+1}|$

One sample Kolmogorov-Smirnov (K-S) statistic.

$$T_{KS} = \max_{1 \leq i \leq n} (\max \left\{ \left| F(R_i) - \frac{i-1}{n} \right|, \left| F(R_i) - \frac{i}{n} \right| \right\}).$$

For large n : $T_{KS} \rightarrow T$

	H_0 True	H_0 False
* Null Hypothesis :	H_0 reject	Type I
	H_0 not reject	Correct
		Type II

α - level of significance.

* Null Distribution & Type I error :

Let us consider our statistic T , how large (or small) T needs to be before we can safely assert that H_0 is false.

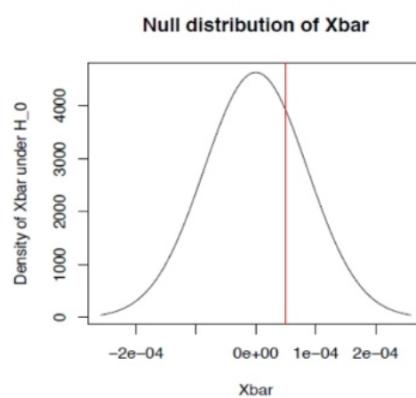
In most cases, we can never be 100% sure that H_0 is false but we can compute T from the observed data and compare it with the sampling distribution of T if H_0 were true. This is called the null distribution of T .

Under H_0 : $\bar{X} \sim N(0, \frac{\sigma_0^2}{n})$ and $\bar{Y} \sim N(0, \frac{\sigma_0^2}{n})$
 $\text{as } (X_i, Y_i) \stackrel{iid}{\sim} N(0, \sigma_0^2 I)$

Based on the observed data $\bar{X} = 0.5 \times 10^{-4}$ or another data ; $\bar{X} = 2.5 \times 10^{-4}$

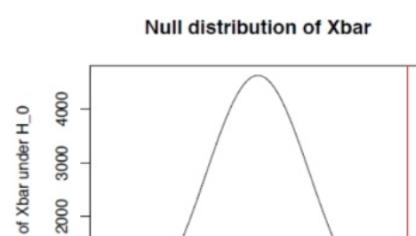
Under H_0 , $\bar{X} \sim N(0, 2.23 \times 10^{-7}/n)$. This normal distribution is the null distribution of \bar{X} .

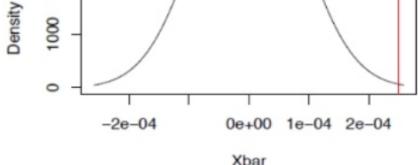
Here's the PDF for the null distribution of \bar{X} , when $n = 30$:



If, for the observed data, $\bar{X} = 0.5 \times 10^{-4}$, this would not provide strong evidence against H_0 . In this case, we might accept H_0 .

Here's the PDF for the null distribution of \bar{X} , when $n = 30$:





If, for the observed data, $\bar{X} = 2.5 \times 10^{-4}$, this would provide strong evidence against H_0 . In this case we might reject H_0 .

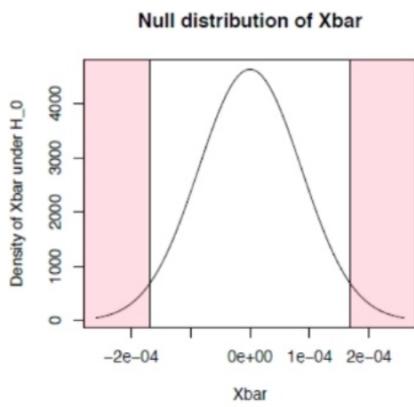
The rejection region is the set of values of T for which we would reject H_0 . The acceptance region is the set of values of T for which we choose to accept H_0 .

$P_{H_0}(\text{reject } H_0) = \alpha$, the value α is the significance level.

Under H_0 , T belongs to the rejection region with probability α then the test is level- α test.

Two sided level- α test reject H when T falls in the shaded region.

(Notation: For a simple null hypothesis H_0 , we write $\mathbb{P}_{H_0}[\mathcal{E}]$ to denote the probability of event \mathcal{E} under H_0 , i.e., the probability of \mathcal{E} if H_0 were true.)



Example 1.6.2. A (two-sided) level- α test might reject H_0 when \bar{X} falls in the above shaded regions. Mathematically, let z_α denote the $1 - \alpha$ quantile, or “upper α point”, of the distribution $\mathcal{N}(0, 1)$. As $\bar{X} \sim \mathcal{N}(0, \sigma^2/n)$ under H_0 (where $\sigma^2 = 2.23 \times 10^{-7}$), the rejection region should be $(-\infty, -\frac{\sigma}{\sqrt{n}} \times z_{\alpha/2}] \cup [\frac{\sigma}{\sqrt{n}} \times z_{\alpha/2}, \infty)$.

Acceptance region : $]-\frac{\sigma}{\sqrt{n}}, z_{\alpha/2}], [\frac{\sigma}{\sqrt{n}}, z_{\alpha/2}[$

Rejection region : $]-\infty, -\frac{\alpha}{\sqrt{n}} z_{\alpha/2}] \cup [\frac{\alpha}{\sqrt{n}} z_{\alpha/2}, +\infty [$

* **Def:** P - Value

The smallest significance level at which your test would be rejecting H_0 .

Let the test be a one sided test that rejects for large values of T if t_{obs} denotes the value of T computed based on the obs data
p-value = $\Pr_{H_0}(T \geq t_{obs})$

Two sided: p-value = $2 \min(\Pr_{H_0}(T \geq t_{obs}), \Pr(T \leq t_{obs}))$

- * P-value is a quantitative measure of the extent to which the data supports (or does not support) H_0 .
- * Failing to reject $H_0 \not\Rightarrow$ there is strong evidence that H_0 is true

- A) The particular test statistic chosen is not good at distinguishing the null hypothesis from the true distribution. Equivalently, the true distribution is not well captured by H_1 that your statistic is capturing.
 - B) You don't have enough data to reject H_0 at the desired significance level. In this case, your study is underpowered.
- * Suppose that $H_0: X_1, \dots, X_n \sim N(8, 1)$ and $T = \bar{X}$ so the exact null distribution is $N(8, 1/n)$.

* Determining the null distribution:

1) Exact Null Distribution

2) Derive an asymptotic approximation, using CLT & continuous mapping theorem (when H_0 is simple, we can obtain null distribution by simulation).

* Deriving asymptotic null distribution :

Let X_1, \dots, X_6 denote the counts of 1, 2, ..., 6 from n rolls of a die, and consider testing if the die is fair.

$$H_0: (X_1, \dots, X_6) \sim \text{Multinomial}(n, (\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}))$$

$$\text{Test statistic : } T = \sum_{i=1}^6 \left(\frac{X_i}{n} - \frac{1}{6} \right)^2$$

$$X_i \sim \text{Bin}(n, \frac{1}{6}) \quad \text{by CLT} \quad \frac{X_i - \frac{n}{6}}{\sqrt{\frac{n \cdot \frac{5}{6} \cdot \frac{1}{6}}{n}}} \sim N(0, 1) \quad \text{and}$$

$$\text{as } n \rightarrow +\infty \quad \frac{X_i - \frac{n}{6}}{\sqrt{n}} = \frac{X_i}{n} - \frac{1}{6} \sim N(0, 1)$$

Hence, $\left(\frac{X_i}{n} - \frac{1}{6} \right) \sim \chi_5^2 \Rightarrow T \sim \chi_5^2$ as (X_1, \dots, X_6) are not jointly independent

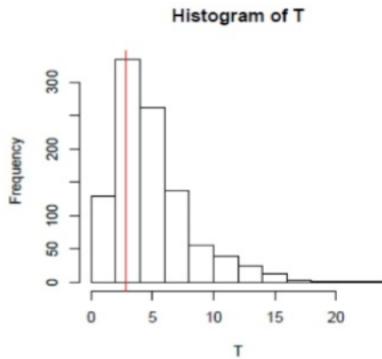
Thus, T approximately distributed as $\frac{1}{6n} \chi_5^2$ distribution

Rejection region is if $T > c$, c is derived from the level of significance. We get c from the $1-\alpha$ quantile of $\frac{1}{6n} \chi_5^2$ so

$$c = \frac{1}{6n} \chi_5^2(\alpha)$$

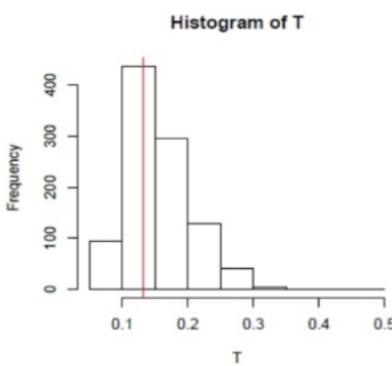
1.11 Using a simulated null distribution

Example 1.11.1. Let T be Pearson's chi-squared statistic for goodness of fit for the values $X_1^2 + Y_1^2, \dots, X_{30}^2 + Y_{30}^2$ from Perrin's experiments. We may simulate the null distribution of T :



This shows the 1000 values of T across 1000 simulations. The observed value $t_{\text{obs}} = 2.83$ for Perrin's real data is in red.

Example 1.11.2. Let T be the K-S statistic for $X_1^2 + Y_1^2, \dots, X_{30}^2 + Y_{30}^2$. We may simulate the null distribution of T :



Page 10

The observed value $t_{\text{obs}} = 0.132$ for Perrin's real data is in red.

We obtain an approximate p -value as the fraction of simulated values of T larger than t_{obs} . (For a two-sided test, we would take either the fraction of simulated values of T larger than t_{obs} or smaller than t_{obs} , and multiply this by 2.)

For Perrin's data, the Pearson chi-squared p -value is 0.754, and the K-S p -value is 0.612. We accept H_0 in both cases, and neither test provides significant evidence against Einstein's theory of Brownian motion.

* Simple alternative and the Neyman-Pearson lemma

- Till now, we discussed a number of ways to construct test statistics for testing a simple null hypothesis, and we showed how to use the null distribution of the statistic to determine the rejection region so as to achieve the desired significance level.
- In this section, our goal is to answer the following question: Which test statistic

In this section, our goal is to answer the following question. Which test statistic should we use?

- The answer depends on the alternative hypothesis that we wish to distinguish from the null.

We will focus on simple H_0 against simple H_1 .

$$\beta = P_{H_0}(\text{Accept } H_0) \quad \text{probability of type 2 error}$$

$$1 - \beta = P_{H_1}(\text{Reject } H_0) \quad \begin{aligned} &\text{probability of correctly rejecting } H_0 \\ &\text{called the power of the test.} \end{aligned}$$

Goal : Maximize the power of the test subject to the significance level of the test under H_0 is at most α

This is a constrained optimization problem

Let x be the data vector, $x = (x_1, \dots, x_n)$ and let $u = (u_1, \dots, u_n)$ be vector of possible values of x .

H_0 : x is distributed with joint PMF $f_0(u) := f_0(u_1, \dots, u_n)$

H_1 : x is distributed with joint PMF $f_1(u) := f_1(u_1, \dots, u_n)$

Let \mathcal{X} be the set of all both values of x under f_0 & f_1 .
We need to specify the rejection region $R \subset \mathcal{X}$ such that we reject H_0 if the observed data belongs to R otherwise we accept H_0 .

$$P_{H_0}(\text{Reject } H_0) = \alpha = \sum_{u \in R} f_0(u)$$

$$P_{H_1}(\text{Reject } H_0) = 1 - \beta = \sum_{u \in R} f_1(u)$$

Formally, we define the optimization problem as choosing the rejection region $R \subset X$ with the goal maximize $\sum_{n \in R} f_0(n)$ subject to $\sum_{n \in R} f_1(n) \leq \alpha$.

We want to know what are the best points n to include in the rejection region R . Hence, R should consist of those points n corresponding to the smallest values of $\frac{f_0(n)}{f_1(n)}$, as these give the smallest increase in type I error for unit increase in the power. These points give most significant evidence against H_0 and in favour of H_1 over H_0 . The statistic $L(x) = \frac{f_0(x)}{f_1(x)}$ is the likelihood ratio statistic and it will reject for small values of $L(x)$.

Example 2.2.1. Consider data X_1, \dots, X_n and the following null and alternative hypotheses:

$$H_0 : X_1, \dots, X_n \stackrel{\text{IID}}{\sim} \mathcal{N}(0, 1)$$

$$H_1 : X_1, \dots, X_n \stackrel{\text{IID}}{\sim} \mathcal{N}(\mu, 1).$$

Here we assume μ is a known, specified value (not equal to 0), so that H_1 is a simple alternative hypothesis. The joint PDF of (X_1, \dots, X_n) under H_0 is

$$f_0(n_1, \dots, n_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-n_i^2/2} = (\alpha\pi)^{-n/2} \cdot e^{-\frac{1}{2} \sum_{i=1}^n n_i^2}$$

$$f_1(n_1, \dots, n_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(n_i-\mu)^2}{2}} = (\alpha\pi)^{-n/2} \cdot e^{-\frac{1}{2} \sum_{i=1}^n (n_i - \mu)^2}$$

$$\begin{aligned} L(n_1, \dots, n_n) &= \frac{f_0(n_1, \dots, n_n)}{f_1(n_1, \dots, n_n)} = \exp \left\{ -\frac{1}{2} \sum_{i=1}^n n_i^2 + \frac{1}{2} \sum_{i=1}^n (n_i - \mu)^2 \right\} \\ &= \exp \left\{ \frac{1}{2} \sum_{i=1}^n \mu^2 - \mu \cdot \sum_{i=1}^n n_i \right\} \\ &= \exp \left\{ \frac{n}{2} \mu^2 - \mu \bar{n} \right\} \end{aligned}$$

$$= \exp \left\{ \frac{n}{2} \gamma (\gamma - 2\bar{u}) \right\}$$

$$= \exp \left\{ \frac{n\gamma^2 - 2\gamma\bar{u}}{2} \right\}$$

If $\gamma > 0$, then $L(n)$ is a decreasing function of \bar{u} . Hence, rejecting for small values of $L(n)$ is equivalent to rejecting large values of \bar{x} .

Rejecting H_0 if $\bar{x} > c$. Under H_0 , $\bar{x} \sim N(0, 1/n)$, $c = \frac{1}{\sqrt{n}} z_\alpha$ where z_α is the upper α point of the $N(0, 1)$.

If $\gamma < 0$, then $L(n)$ is an increasing function of \bar{u} . Hence, rejecting for small values of $L(n)$ is equivalent to rejecting small values of \bar{x} .

Rejecting H_0 if $\bar{x} < c$. Under H_0 , $\bar{x} \sim N(0, 1/n)$, $c = \frac{-1}{\sqrt{n}} z_\alpha$ where $-z_\alpha$ is the upper $1-\alpha$ point of the $N(0, 1)$. $\bar{x} < \frac{-z_\alpha}{\sqrt{n}}$

If we change it to $H_1: x_1, \dots, x_n \sim N(18, 1)$ and $H_1: x_1, \dots, x_n \sim N(10^3, 1)$ as we are not considering the value of α then we only consider the case $\gamma > 0$.

In this case as $\gamma \neq 0$ so we have two powerful tests one of each side so we do not have one uniformly powerful test $\forall n$

* Lmmer : Neyman Pearson

Let H_0 and H_1 be simple hypotheses, for a constant $c > 0$, suppose that the likelihood ratio test rejects H_0 when $L(n) < c$ has significance level α . Then, for any other test of H_0 with at most significance level α , its power

against H_1 is at most the power of the likelihood ratio test.

* Proof :

Consider the discrete case and let $R = \{n \mid L(n) < L_0\}$ be the rejection region of the Likelihood Ratio Test. Note that among all the subsets of \mathcal{X} , R maximize the quantity

$$\sum_{n \in R} (c f_1(n) - f_0(n)) \text{ because } c f_1(n) - f_0(n) > 0 \quad \forall n \notin R \text{ and} \\ c f_1(n) - f_0(n) \leq 0 \quad \forall n \in R.$$

Let R' be another rejection region. Then,

$$\sum_{n \in R'} (c f_1(n) - f_0(n)) \leq \sum_{n \in R} (c f_1(n) - f_0(n))$$

$$\Rightarrow c \left(\sum_{n \in R} f_1(n) - \sum_{n \in R'} f_1(n) \right) \geq \sum_{n \in R} f_0(n) - \sum_{n \in R'} f_0(n) \\ = \alpha - \sum_{n \in R'} f_0(n) \leq \alpha \\ \geq 0$$

Hence, $\sum_{n \in R} f_1(n) \geq \sum_{n \in R'} f_1(n)$ so we deduce the power of the LRT is at least that of the other test. Similarly, for the continuous case.

* Uniformly Most Powerful Test :

The most powerful test against the alternative $H_1: x_1, \dots, x_n \sim N(\mu, 1)$ is the same for $\mu \geq 0$ (reject if $\bar{x} > \frac{1}{\sqrt{n}} z_\alpha$) and neither the test nor the test statistic nor the rejection region depends on the specific value of μ . This means that, in fact, this test is uniformly most powerful (UMP) against the one-sided composite hypothesis.

$H_0: X_i \sim N(\mu, 1) \forall i=1, \dots, n$ where $\mu \geq 0$ or $X_i \sim N(\mu, 1)$ where $\mu < 0$

On the other hand, the MP test is different for $\mu \geq 0$ vs. $\mu < 0$. one rejects for large positive values of \bar{X} and the other rejects for large negative values of \bar{X} . Hence, a single MP test does not exist for the two sided hypothesis $H_0: X_i \sim N(\mu, 1) \forall i=1, \dots, n, \mu \neq 0$



* Example :

Find the MP test at level of significance α to test whether the data $X_1, \dots, X_n \sim \text{Ber}(p)$ where $H_0: p = \frac{1}{2}$ $H_1: p \neq \frac{1}{2}$

$$L(X) = \frac{f_0(x_1, \dots, x_n)}{f_1(x_1, \dots, x_n)} = \frac{\prod_{i=1}^n (p)^{x_i} (1-p)^{1-x_i}}{\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}} = \frac{1}{\alpha^n} \frac{1}{p^{\sum x_i} (1-p)^{n-\sum x_i}}$$

$$L(X) = \frac{1}{\alpha^n (1-p)^n} \left(\frac{1-p}{p} \right)^{\sum_{i=1}^n x_i}$$

$p \leq \frac{1}{2}$: Hence, $L(X)$ is an increasing function of $S = \sum_{i=1}^n x_i$ so we reject for small values of S . By NP lemma, the most powerful test rejects H_0 when $S \leq c$ for some c .

$$H_0: X_1, \dots, X_n \sim \text{Bernoulli}(p), p \leq \frac{1}{2}$$

$S \sim \text{Binomial}(n, \frac{1}{2})$ under H_0 so c is the α quantile of $\text{Binomial}(n, \frac{1}{2})$.

$P > \frac{1}{2}$: Hence, $L(X)$ is a decreasing function of $S = \sum_{i=1}^n X_i$ so we reject for large values of S . By NP lemma, the most powerful test rejects H_0 when $S > c$ for some c .

$$H_1 : X_1, \dots, X_n \sim \text{Bernoulli}(p), \quad P > \frac{1}{2}$$

$S \sim \text{Binomial}(n, \frac{1}{2})$ under H_0 so c is the $(1-\alpha)$ quantile of $\text{Binomial}(n, \frac{1}{2})$.

This test is the same for all $p < \frac{1}{2}$ ($p > \frac{1}{2}$) and it is in fact the uniformly most powerful test against the one sided composite hypothesis $H_1 : X_1, \dots, X_n \sim \text{Ber}(p)$ ($p < \frac{1}{2}$, $p > \frac{1}{2}$).

* Remark K :

We have glanced over a detail which is that when the distribution of $L(X)$ is discrete under H_0 , it might not be possible to choose c ; the significance level exactly α .

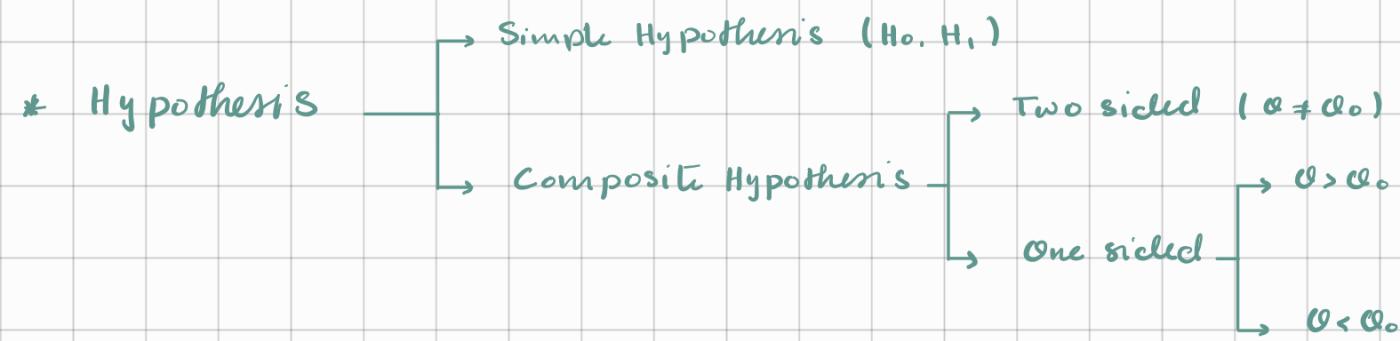
Example: $S \sim \text{Bin}(20, \frac{1}{2})$

$$P(S \geq 15) = 0.021 \quad P(S \geq 14) = 0.088$$

If we reject H_0 when $S \geq 14$, we don't reach the significance level α . If we reject H_0 when $S \geq 15$, we are too conservative.

Theoretically, we do a randomized test. Always reject H_0 when $S \geq 15$, we accept H_0 when $S \leq 13$ and reject H_0 with

at certain probability when $S=14$ where this probability is chosen to make the significance level is exactly α . By NP lemma, this is the most powerful test among all randomized tests. In practice, it might not be acceptable to use randomized test, so we might take the more conservative option when $S \geq 15$.



* Example : Suppose that there are 80 students in a class taking MATH 350 course. A diagnostic example is administered at the start of the quarter, and a comparable exam is administered at the end of the quarter. Did the MATH 350 course improve students' knowledge of statistics?

Let X_i be the difference of scores of the two tests for the i^{th} student. Assume $X_1, \dots, X_{80} \sim N(\mu, \sigma^2)$.

We might formulate the hypothesis that $H_0: \mu = 0$ vs. $H_1: \mu > 0$. Both hypotheses are composite as σ^2 is unknown.

If we cannot make the normality assumption then we might assume that X_1, \dots, X_{80} are iid with some pdf f and then we test.

H_0 : f is symmetric around zero

H_1 : f is symmetric around μ

We might drop the symmetry assumption of the test. Then,

H_0 : f has median zero
 H_1 : f has median $\neq 0$

Selection of the test : 1) Some prior knowledge on the distribution of the scores.
2) Visual Implementation

It's hard to make two exams that are equally difficult. What if the second exam happened to be easier?

To address this we add a control group. We give 100 other students who are not taking MATH 350 course the same two exams at the start & at the end of the quarter. Let y_i be the difference in scores for the student i for the control group. If we believe that $y_1, y_2, \dots, y_{100} \stackrel{iid}{\sim} N(\mu_y, \sigma^2)$ and $x_1, \dots, x_n \stackrel{iid}{\sim} N(\mu_x, \sigma^2)$.

To do this, we might formulate the test as $H_0: \mu_x = \mu_y$
 $H_1: \mu_x > \mu_y$

If we cannot have the normality assumption then we might test
 $H_0: f = g$ where f is the PDF of the x_i 's
 $H_1: f$ stochastically dominates g and g is the PDF of the y_i 's.

We can rewrite the alternative hypothesis as $P(x \geq n) \geq P(y \geq y)$ $\forall n \in \mathbb{R}$.

We have to inspect the data again to check the normality assumption.

* Computing the value of α :

$$x_1, \dots, x_n \sim N(\mu, \sigma^2)$$

$$\text{Test statistic : } \bar{x}$$

Reject H_0 : $\bar{x} > c$

H_0 : $\mu = 80$

H_1 : $\mu > 80$

Under H_0 : $\bar{x} \sim N(80, \frac{\sigma^2}{n})$

$$\alpha = P_{H_0}(\bar{x} > c), \alpha = 0.05$$

$$= P\left(\frac{\sqrt{n}(\bar{x}-80)}{\sigma_0} > \frac{\sqrt{n}(c-80)}{\sigma_0}\right)$$

$$= P(Z \geq c^*)$$

$$0.95 = P(Z \leq c^*)$$

$$\Rightarrow 1.65 = \frac{\sqrt{n}(c-80)}{\sigma_0}, \sigma_0 = 10$$

$$n = 100$$

$$\Rightarrow c = 81.65$$

Rejecting H_0 when $\bar{x} > 81.65$ with 0.05 significance level.

Suppose that $x_1, \dots, x_n \sim N(\mu, \sigma^2)$; $H_0: \mu = 80$ $H_1: \mu > 80$ where σ^2 is unknown. Composite H_0 & H_1 .

What is the distribution of the test statistic \bar{x} ? Calculate Type I error.

When testing a composite H_0 against a composite H_1 . When testing a composite H_1 , there is a probability of Type I error associated to each of the data distribution $P \in H_0$ (the prob of rejecting H_0 if the true distribution is P) and a probability of type II error associated to each of the data distribution $P \in H_1$ (prob of accepting H_0 when P is the true distribution). A test has significance level α if the maximum probability of Type I error for any $P \in H_0$ is α . This means that to design a level α -test, we need to control the prob. of type I error for every $P \in H_0$. This is difficult!

* One sample t-test :

Assume $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$ for unknown μ and σ^2 . Consider testing :

$$H_0 : \mu = 0$$

$$H_1 : \mu > 0$$

σ^2 is unknown, a natural idea is to estimate σ^2 by the sample variance : $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ and consider the statistic :

$$T = \frac{\sqrt{n}(\bar{X} - \mu_0)}{S} = \frac{\sqrt{n}(\bar{X}/\sigma)}{S/\sigma} \sim \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2 / n-1}} = t_{n-1}, \text{ if } X_i^2 \text{ are independent}$$

The distribution of t does not depend on σ , so it is same for any μ_0 .

* Def : t-distribution

If $Z \sim N(0, 1)$ and $V \sim \chi^2_{n-1}$; Z and V are independent then the distribution of $\frac{Z}{\sqrt{V/n}}$ is called the t-distribution with $n-1$ df denoted by t_{n-1} .

Under H_0 , $T \sim t_{n-1}$. Let $t_{n-1, \alpha}$ denote the upper α point of the distribution t_{n-1} , then the test that rejects for $T > t_{n-1, \alpha}$ is called the one sample t-test.

The one sample t-test is often used in paired two sample setting; the previous example. There we actually have two paired samples before and after the test per student. We perform the test by first considering the differences of the

scores, such settings are called paired two sample test although we perform a one sample t-test.

* Two sample t-test :

Let two independent samples $x_1, \dots, x_n \stackrel{iid}{\sim} N(\mu_x, \sigma^2)$ and $y_1, \dots, y_m \stackrel{iid}{\sim} N(\mu_y, \sigma^2)$. Suppose that they have a common variance.

$$H_0: \mu_x - \mu_y = 0$$

$$H_1: \mu_x - \mu_y > 0$$

A natural idea is to reject H_0 for large values of $\bar{x} - \bar{y}$. $\bar{x} \sim N(\mu_x, \frac{\sigma^2}{n})$
 $\bar{y} \sim N(\mu_y, \frac{\sigma^2}{m}) \Rightarrow \bar{x} - \bar{y} \sim N(\mu_x - \mu_y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m})$

$$\text{Under } H_0: T \sim N(0, \sigma^2 (\frac{1}{n} + \frac{1}{m})) \Rightarrow \frac{\bar{x} - \bar{y}}{\sqrt{\sigma^2 (\frac{1}{n} + \frac{1}{m})}} \sim N(0, 1)$$

- If σ^2 is known : $\frac{\bar{x} - \bar{y}}{\sqrt{\sigma^2 (\frac{1}{n} + \frac{1}{m})}} \rightarrow z_\alpha$ then we reject H_0

- If σ^2 is unknown : $S_p^2 = \frac{1}{m+n-2} \left(\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^m (y_i - \bar{y})^2 \right)$
is called the pooled sample variance and it is used to estimate σ^2 .

$$\frac{S_p^2}{\sigma^2} \sim \frac{1}{m+n-2} \chi^2_{m+n-2}$$

$$T = \frac{(\bar{x} - \bar{y}) / \sqrt{\sigma^2 (\frac{1}{n} + \frac{1}{m})}}{\sqrt{S_p^2 / \sigma^2}} = \frac{\bar{x} - \bar{y}}{\sqrt{S_p^2 (\frac{1}{n} + \frac{1}{m})}} \sim t_{n+m-2}$$

$$P_{H_0}(T > c) = \alpha \Rightarrow P(T > t_{n+m-2, \alpha}) = \alpha$$

$$\Rightarrow c = t_{n+m-2, \alpha}$$

We reject H_0 when $T > t_{n+m-2, \alpha}$ (upper α ($1-\alpha$) point of t_{n+m-2}) for the α sample t-test.

* **Remark :** The assumption of common variance σ^2 for the samples is often violated in practice. We can assume that $x_1, \dots, x_n \sim N(\mu_x, \sigma_x^2)$ and $y_1, \dots, y_m \sim N(\mu_y, \sigma_y^2)$ for $\sigma_x^2 \neq \sigma_y^2$.

$\text{Var}(\bar{x} - \bar{y}) = \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}$, we may estimate this by $\frac{s_x^2}{n} + \frac{s_y^2}{m}$ where s_x^2, s_y^2 are the sample variances.

Then, we may use the test statistic :

$$\text{Welch} = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}$$

Welch under H_0 is not exactly t-distribution but it was shown by Welch to be close to a t-distribution with

$$\frac{(s_x^2/n + s_y^2/m)^2}{\frac{(s_x^2/n)^2}{n-1} + \frac{(s_y^2/m)^2}{m-1}} \text{ d.f.}$$

The test that rejects H_0 when Welch exceeds the upper α point of this distribution is called the Welch t-test or the unequal variances t-test.

* Testing of Proportions :

lets assume you roll a 4-faced die with values {1, 2, 3, 4} 100 times and the value "4" appears 29 times, does this clearly suggest your die is unbiased.

If the die is unbiased we can expect the proportion of "4" to be 0.25. Let the proportion of "4" be P .

$$H_0: P_0 = 0.25 \quad \text{vs} \quad H_1: P_1 \neq 0.25$$

We observe the proportion of 4 to be $\hat{P} = 0.29$.

Let Y be the number of times "4" appears Under H_0 :

$$\mathbb{E}(Y) = n\hat{P} = nP_0$$

$$\text{Var}(Y) = n\hat{P}(1-P_0) = nP_0(1-P_0)$$

$$\mathbb{E}(\hat{P}) = \mathbb{E}\left(\frac{Y}{n}\right) = \frac{1}{n} \mathbb{E}(Y) = P_0$$

$$\text{Var}(\hat{P}) = \text{Var}\left(\frac{Y}{n}\right) = \frac{1}{n^2} \text{Var}(Y) = \frac{P_0(1-P_0)}{n}$$

By CLT: $\hat{P} \sim N(P_0, \frac{P_0(1-P_0)}{n})$ as $n \rightarrow +\infty$

$$Z = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}} \rightarrow \text{observed} \quad \text{Critical: } Z_{\alpha/2}$$

Acceptance Region: $[-Z_{\alpha/2}, Z_{\alpha/2}]$

* Testing of Proportion:

Times magazine reported the results of a telephone poll of 800 adults. The question was "Should America increase the federal tax on cigarettes to pay for health care reforms?". The results are as follows:

	Non-Smokers (1)	Smokers (2)
--	-----------------	-------------

n_1 : # of participants	605	195
y_1 # of "yes" response	351	41

Can you conclude that the 2 groups of individuals significantly differ in their opinion?

Let p_1, p_2 denote the proportion of individuals responding "yes" in the two subgroups.

$$\hat{p}_1 = \frac{y_1}{n_1} = \frac{351}{605} = 0.58$$

$$H_0: p_1 = p_2$$

$$H_1: p_1 \neq p_2$$

$$\hat{p}_2 = \frac{y_2}{n_2} = \frac{41}{195} = 0.21$$

$$\text{By CLT: } \hat{p}_1 \sim N\left(p_1, \frac{p_1(1-p_1)}{n_1}\right)$$

$$\hat{p}_2 \sim N\left(p_2, \frac{p_2(1-p_2)}{n_2}\right)$$

Pooled Estimate of p

$$\hat{p}_1 - \hat{p}_2 \sim N\left(0, p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)\right), \quad p = \frac{y_1 + y_2}{n_1 + n_2}$$

$$z_{\text{obs}} = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} > z_{\alpha/2} \text{ or } z_{\text{obs}} < -z_{\alpha/2} \text{ reject.}$$

* The theoretical work on the hypothesis testing of variance is primarily based on fitting a Normal distribution to the data & then the use of sampling distribution.

* One Sample Test :

Given $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. We aim to test whether

$H_0: \sigma^2 \leq \sigma_0^2$ against the alternatives

$H_1: \sigma^2 \neq \sigma_0^2, \sigma^2 > \sigma_0^2, \sigma^2 < \sigma_0^2$

$$S^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Under $H_0: \frac{(n-1)S^2}{\sigma_0^2} \sim \chi_{n-1}^2$

Rejection Region : $\sigma^2 \neq \sigma_0^2$

$$\chi^2 > \chi_{\alpha/2, n-1}^2 \text{ or } \chi^2 < \chi_{1-\alpha/2, n-1}^2$$

$$\sigma^2 > \sigma_0^2$$

$$\chi^2 > \chi_{\alpha, n-1}^2$$

$$\sigma^2 < \sigma_0^2$$

$$\chi^2 < \chi_{1-\alpha, n-1}^2$$

- * A manufacturer of helmet safety belts is concerned about the mean and variation of the force its helmets transmits to wearers when subjected to external force. They designed the helmets so that the mean force transmitted is 800 pounds with S.D to be less than 40 pounds. Tests on 40 helmets were run and $\bar{X} = 825$ pounds and $S = 48.5$ pounds. Do the data provide sufficient evidence that the population SD exceeds 40 pounds. Do the data provide sufficient evidence at $\alpha = 0.05$ to conclude that $SD \neq 40$ pounds.

$$\chi_{0.05, 39}^2 = 54.572$$

$$\chi_{1-0.025, 39}^2 = 83.654$$

$$\chi_{0.025, 39}^2 = 58.120$$

$$a) H_0: \sigma = 40 \quad H_1: \sigma > 40$$

$$\chi^2 = \frac{(n-1) S^2}{\sigma^2} = \frac{(40-1) (48.5)^2}{(40)^2} = 57.34 > \chi^2_{0.05, 39}$$

Hence, the standard deviation is greater than 40.

$$b) H_0: \sigma = 40 \quad H_1: \sigma \neq 40$$

We have that $\chi^2_{1-\alpha/2, 39} < \chi^2 < \chi^2_{\alpha/2, 39}$ so we cannot reject the null hypothesis. It does not significantly differ from 40.

* Two Sample Test :

Let $N(\mu_x, \sigma_x^2)$ and $N(\mu_y, \sigma_y^2)$ be two independent populations. x_1, \dots, x_n is a random sample from $N(\mu_x, \sigma_x^2)$ and y_1, \dots, y_m is a random sample from $N(\mu_y, \sigma_y^2)$ ($n \neq m$)

$$H_0: \sigma_x^2 = \sigma_y^2$$

$$H_1: \sigma_x^2 < \sigma_y^2, \quad \sigma_x^2 > \sigma_y^2, \quad \sigma_x^2 \neq \sigma_y^2$$

$$\frac{S_x^2(n-1)}{\sigma_x^2} \sim \chi^2_{n-1}, \quad \frac{S_y^2(m-1)}{\sigma_y^2} \sim \chi^2_{m-1}$$

$$\frac{\chi^2_x}{\chi^2_y} = \frac{S_x^2(n-1) \sigma_y^2}{S_y^2(m-1) \sigma_x^2} \sim \chi^2_{m-1}$$

$$\text{Under } H_0: \frac{S_x^2}{S_y^2} = \frac{\chi^2_x / (n-1)}{\chi^2_y / (m-1)} \sim F_{n-1, m-1}$$

Rejection Region : $\sigma_x^2 > \sigma_y^2$

$$F > F_{\alpha, n-1, m-1}$$

$$\sigma_x^2 < \sigma_y^2$$

$$F < F_{1-\alpha, n-1, m-1} = \frac{1}{F_{\alpha, m-1}}$$

$$\sigma_x^2 \neq \sigma_y^2$$

$$F > F_{\alpha/2, n-1, m-1} \text{ or } F < F_{\alpha/2, n-1, m-1}$$

- * A scientist is interested in whether male and female have driving behaviours. The question she framed was as follows "Is the mean fastest speed driven by a male student different than the mean fastest speed driven by the female students. The psychologist conducted a survey of 34 male students and 29 female students. The statistics she received :

	Male	Female
Mean	105.5	90.9
Sd	20.1	12.2

Is there sufficient evidence at $\alpha = 0.05$ level to conclude that the variance of the fastest speed driven by male student differs from the fastest speed driven by female students.

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

1 : Male

2 : Female

$$\text{Under } H_0: \frac{S_1^2 (n-1)}{S_2^2 (m-1)} = \frac{(20.1)^2}{(12.2)^2} = 2.714 = F$$

$$F_{0.025, 33, 28} = 2.089$$

$$F_{1-0.025, 33, 28} = \frac{1}{2.089}$$

$$F > F_{0.025, 33, 28} \text{ so we reject}$$

the null so the fastest mean differs

between male and female.

* Question 1:

Let $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(p)$ denote n tosses of a biased coin and assume that $\hat{p} = \bar{x}$. We will explore two ways to construct a 95% confidence interval both based on CLT. $\sqrt{n}(\hat{p} - p) \rightarrow N(0, p(1-p))$

- a) Use a plugin estimate $\hat{p}(1-\hat{p})$ for the variance $p(1-p)$ to obtain a 95% CI for p . Use the Wald confidence interval.

By CLT $E(p) = \hat{p}$ and $\text{Var}(p) = \hat{p}(1-\hat{p})$

$$\frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \sim N(0, 1) \Rightarrow P_p(-z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \leq z_{\alpha/2}) = 0.95$$

$$\Rightarrow p \in \left[\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

- b) We can use the CLT which implies for large n

$$P_p(-\sqrt{p(1-p)} z_{\alpha/2} \leq p \leq \sqrt{p(1-p)} z_{\alpha/2}) = 1 - \alpha$$

Solve the equation $\sqrt{n}(\hat{p} - p) = \sqrt{p(1-p)} z_{\alpha/2}$ for p in terms of \hat{p} and also solve $\sqrt{n}(\hat{p} - p) = -\sqrt{p(1-p)} z_{\alpha/2}$ for p in terms of \hat{p} to obtain a 95% CI for p .

$$n(\hat{p} - p)^2 \leq p(1-p) z_{\alpha/2}^2$$

$$n\hat{p}^2 - 2np\hat{p} + np^2 \leq p z_{\alpha/2}^2 - p^2 z_{\alpha/2}^2$$

$$(n + z_{\alpha/2}^2)p^2 - (2n\hat{p} + z_{\alpha/2}^2)p + n\hat{p}^2 \leq 0$$

p is between the two real roots of this equation

$$\Delta = (2n\hat{p} + z_{\alpha/2}^2) - 4(n + z_{\alpha/2}^2)n\hat{p}^2$$

$$\Rightarrow p \in \left[\frac{(2n\hat{p} + z_{\alpha/2}^2) \pm \sqrt{\Delta}}{2(n + z_{\alpha/2}^2)} \right]$$

$$\alpha(n + \alpha_2)$$

c) Perform a simulation study to determine the true coverage at the CIs at parts a and b for the combinations:

$$n = (10, 40, 100)$$

$$p = (0.1, 0.3, 0.5)$$

for each combination perform at least $B = 100,000$ simulations in each simulation you may simulate $\hat{p} \sim \text{Bin}(n, p)$ instead of simulating x_1, \dots, x_n . Report the simulated coverage of the both CIs.

* Question 2:

In a classic genetic study, it was recorded that hospital records gender ratios. The following table shows # of male children out of 6115 families each having 12 children.

# of Male children	# families
0	7
1	45
2	181
3	478
4	829
5	1112
6	1343
7	1033
8	670
9	286
10	104
11	24
12	3

Let x_1, \dots, x_{6115} denote the # of male children.

a) Suggest two reasonable tests for testing $H_0: x_1, \dots, x_{6115} \stackrel{iid}{\sim} \text{Bin}(12, \frac{1}{2})$

$$T_2 = \sum_{i=0}^{12} \frac{(O_i - E_i)^2}{E_i} = \sum_{i=0}^{12} \frac{(O_i - 6115 \binom{12}{i} (\frac{1}{2})^{12})^2}{6115 \binom{12}{i} (\frac{1}{2})^{12}}$$

$$T_1 = \frac{\# \text{ of males} \times \text{frequency}}{\sum \text{frequency}} \quad \text{mean}(\# \text{ of males})$$

b) Perform a simulation to simulate the null hypothesis $B = 10^5$. Plot the histograms of T_1 and T_2 . Using the simulated values compute the approximate p-values of the tests.

c) In this example, why won't the null hypothesis not work.

No, because there are biological and sociological factors affecting this ratio.

* Question :

Consider the problem of test $x_1, \dots, x_n \sim N(0, 1)$ and $H_0: x_1, \dots, x_n \sim N(\mu, 1)$ where $\mu > 0$. We can use the LRT and the t-test.

a) For $n = 100$ verify numerically that these tests have significance level $\alpha = 0.05$ by performing 10^5 simulations and draw a sample from $N(0, 1)$ compute the test statistics relevant to the LRT (reject $H_0: \bar{x} > z_{0.05}$) and t-test (reject $H_0: \frac{\sqrt{n}\bar{x}}{s} > t_{n-1, 0.05}$) and record whether each test accepts /

rejects H_0 . Record the fraction of simulations that rejects 0.05 and check that it is close to 0.08. set.seed(1)

Add to them W^+ , S. Wilcoxon signed rank test rejects H_0 when $W^+ > \frac{n(n+1)}{4} + \sqrt{\frac{n(n+1)(n+2)}{24}} Z_{0.05}$.

For the sign test, reject H_0 if $S > \frac{n}{2} + \sqrt{\frac{n}{4}} Z_{0.05}$ where S is the # of +ve values in x_1, \dots, x_n .

- b) For $n = 100$, numerically compute the powers of the test against the alternative H_1 for values of $\mu = (0.1 : 0.4)$ do this by performing 10^5 simulations by drawing $x \sim N(\mu, 1)$. Report your findings in a table or visual form.
- c) How do the powers of the 4 test compare when testing against H_0 . You should always use the testing procedure that makes the few dist. assumptions because in practise we never know whether $\sigma^2 = 1$ or the dist is truly normal. Comment on this statement. John A Rice says it has been shown that the Wilcoxon signed test is nearly as powerful as t-test thus use it for small n . Do your simulated results support this statement.

* let's start with one sample setting when $x_1, \dots, x_n \stackrel{iid}{\sim} f$, where we drop the normality assumption and we wish to test.

H_0 : f is symmetric around 0

H_1 : f is symmetric around μ for some $\mu > 0$

Because the shape of f is arbitrary under H_0 , the distribution of the t-statistic is no longer the same under every data dist.

H_0 in particular, it can be very far from H_0 . If n is moderately small and δ is heavy-tailed. We consider instead the signed rank statistic W_+ defined in the following way:

- 1) Sort $|x_1|, |x_2|, \dots, |x_n|$ in an increasing order. Assign the smallest value a rank of 1, the next 2, and the largest valued a rank n .
- 2) Define W_+ as the sum of the ranks corresponding to only the +ve values of x_1, x_2, \dots, x_n .

* Example : Wilcoxon Sign Rank Test

$$x_i : 5, -2, 3, 7$$

$$|x_i| : 5, 2, 3, 7$$

$$\text{Rank} : 3, 1, 2, 4$$

$$\begin{aligned} W_+ &= \sum_i \text{Rank of } |x_i|, x_i > 0 \\ &= 3 + 2 + 4 \\ &= 9 \end{aligned}$$

We expect W_+ to be larger under H_1 than under H_0 because higher rank obs are more likely to be +ve under H_1 . The test that rejects for large W_+ has the same distribution under every H_0 and provides a method for determining the null distribution and rejection threshold for W_+ when n is large. When n is small, we can determine the exact null dist. of W_+ by computing W_+ for all 2^n possible combinations of + and - signs of ranked data.

* **Theorem** : The dist of W_+ is the same for every pdf that is symmetric around 0 for large n , this distribution is asymptotically $N\left(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24}\right)$ i.e

$$\frac{24}{n^2} \left(W_+ - \frac{n(n+1)}{4} \right) \xrightarrow{d} N(0, 1) \text{ as } n \rightarrow +\infty$$

* Proof :

We will show that the distribution of W_+ is the same for every δ and that $E(W_+) = \frac{n(n+1)}{4}$ and $\text{Var}(W_+) = \frac{n(n+1)(2n+1)}{24}$

Let $f_{(x_1, \dots, x_n)} = \prod_{i=1}^n f(x_i)$ be the joint pdf of the data. By symmetry of f about 0, $f(\pm x_1, \dots, \pm x_n)$ is the same for each of the 2^n combinations of +/ - signs. This implies that conditional on $|x_1|, \dots, |x_n|$, the signs of x_1, \dots, x_n are independent and each equal to + or - with prob $\frac{1}{2}$.

Let $I_K = 1$, if the value of the rank K is +ve and $I_K = 0$ if the value is -ve. Thus, $I_1, \dots, I_n \stackrel{iid}{\sim} \text{Ber}(\frac{1}{2})$ for any pdf f that is symmetric about zero. The signed rank statistic is $W_+ = \sum_{k=1}^n K I_k$ since I_1, \dots, I_n have the same distribution under any symmetric pdf about 0, the distribution of W_+ is the same for all such pdf's f .

$$E(W_+) = \sum_{k=1}^n K E(I_k) = \sum_{k=1}^n \frac{K}{2} = \frac{1}{2} \cdot \frac{n(n+1)}{2} = \frac{n(n+1)}{4}$$

$$\text{Var}(W_+) = \sum_{k=1}^n K^2 \text{Var}(I_k) = \frac{1}{4} \cdot \frac{n(n+1)(2n+1)}{6} = \frac{n(n+1)(2n+1)}{24}$$

To explain why W_+ is asymptotically normally distributed, we define the empirical cdf of $|x_1|, |x_2|, \dots, |x_n|$ by :

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{|x_i| \leq t\}$$

$F_n(t)$ is the fraction of values of $|x_i|$ that are at most t . The rank of $|x_i|$ can be computed as $n F(|x_i|)$ so $W_+ = \sum_{i=1}^n n F_n(|x_i|) \mathbb{I}|x_i| > 0$

when n is large, we may show that $F_n(t)$ is with high probability close to the true CDF $F(t)$ of $|x_i|$ for every $t \in \mathbb{R}$, and hence the difference between W_+ and $\tilde{W}_+ = \sum_{i=1}^n n F(|x_i|) \mathbb{I}(x_i > 0)$ is negligible. \tilde{W}_+ is the sum of i.i.d. RVs $Y_i = n F(|x_i|) \mathbb{I}(x_i > 0)$ and hence asymptotically normal by CLT.

$$\text{By CLT : } \frac{W_+ - \left(\frac{n(n+1)}{4}\right)}{\sqrt{\frac{n(n+1)(n+2)}{24}}} \xrightarrow{d} N(0, 1) \text{ as } n \rightarrow +\infty$$

* Rank Sum Test :

The idea of converting observed data values to just ranks so as to deal with heavily skewed data and deviations from normality, can be extended to the test sample setting. Consider two indep. samples $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f$ and $y_1, y_2, \dots, y_m \stackrel{iid}{\sim} g$ where f and g are any two arbitrarily pdfs and we test :

$$H_0 : f = g$$

$$H_1 : f \text{ stochastically dominates } g$$

In otherwords, the alternative hypothesis suggests that values drawn from f tend to be larger than values from g . The rank sum test statistic is defined as :

- 1) Consider the pooled sample of all obs. $x_1, \dots, x_n, y_1, \dots, y_m$. Sort them in increasing order and assign the rank 1 to the smallest value 1 and rank $n+m$ for the largest.
- 2) Define the rank sum test statistic T_y as the sum of ranks corresponding to only y_i values i.e. values from the second sample.

x_i	y_i	Values :	1	2	4	7	8	10
ex.		Ranks :	1	2	3	4	5	6
		T_y :	1 + 4 + 5 = 10					
10	8							

We expect T_y to be smaller under H_1 than under H_0 because H_1 , the values of y_i tend to have smaller ranks. The test that rejects for small values of T_y is called as Wilcoxon rank sum test or alternatively as Mann - Whitney Wilcoxon test. If the alternative $H_1 : f \neq g$ then we would reject for both large & small values of T_y .

The following theorem states that T_y has the same distribution for every $P \in H_0$ and provides a method for determining the null hypothesis and rejection threshold when n and m are both large.

* Theorem :

The distribution of T_y is the same under any pdf $f = g$ for large n and m this dist is approximately $N\left(\frac{n(m+n)}{2}, \frac{nm(n+m+1)}{12}\right)$.

If $f = g$, then each ordering of $x_1, \dots, x_n, y_1, \dots, y_m$ is equally likely. Since T_y depends only on this ordering, its dist must be the same under every pdf $f = g$.

If $I_k = 1$, if the k^{th} largest value of $x_1, \dots, x_n, y_1, \dots, y_m$ belongs to the second sample and $I_k = 0$ otherwise $I_k = 0$. Then, $T_y = \sum_{k=1}^{m+n} k I_k$. Under H_0 , I_k indicates whether the k^{th} individual is selected in a simple of random sample of size m (without replacement) from a population of distribution of size $m+n$.

$$E(T_Y) = \frac{mn(m+n+1)}{12}$$

Hence the $I_{k\ell}$'s are not independent

$$\text{Var}(T_Y) = \frac{mn(m+n+1)}{12}$$

* Permutation & Randomized Test :

The main idea behind the one-sample Wilcoxon signed rank test & the two-sample ranked sum test is to exploit the symmetry under H_0 . For the sign-rank test, the symmetry is that it is equally likely to observe $\pm x_1, \dots, \pm x_n$ for each of the 2^n combinations of \pm -signs. For the ranks-sum test, the symmetry is that it is equally likely to observe each of the $(m+n)!$ permutations of the paired sample $x_1, \dots, x_n, y_1, \dots, y_m$. In fact, this idea of exploiting the symmetry provides an alternative simulation based method of obtaining the null dist. for any test statistic T .

* Example :

Consider $x_1, \dots, x_n, y_1, \dots, y_m$ and any test statistic depending on x_i 's and y_j 's ($T = \bar{x} - \bar{y}$ as e.g.). For a null hypothesis H_0 which specifies that all data from both samples are iid from a common dist.

$$H_0: x_1, \dots, x_n, y_1, \dots, y_m \stackrel{iid}{\sim} f \text{ for any unknown pdf } f.$$

The permutation null dist. of T is the dist. of $T(x_1^*, \dots, x_n^*, y_1^*, \dots, y_m^*)$ when we fix the observed values of $x_1, \dots, x_n, y_1, \dots, y_m$ and let x_1^*, \dots, y_m^* be a permutation of x_1, \dots, y_m chosen at random.

uniformly at random from a set of all $(m+n)!$ possible permutations. Under H_0 , each of these $(m+n)!$ possible values of T is equally likely to be observed. To perform a test that rejects for the large values of T , we may use the following procedure:

- 1) Randomly permute the paired data $B = 10^4$ times and compute the value of T each time.
- 2) Compute an approximate p-value as a fraction of B simulations where we obtained a value of T larger than T_{observed} , the value for the original (unpermuted) data. Reject at level α if the p-value is at most α .

For a two-sided test that rejects for both large and small values of t , we can compute the p-value by taking the frac of sim where $t > t_{\text{obs}}$ or the frac of sim where $t < t_{\text{obs}}$ and multiply by 2. This is called a permutation based on t and an example of unconditional test because we are looking at the conditional data under H_0 given the set of their values. The utility of this idea may be applied to test statistic t where we don't understand its unconditional distribution under H_0 and where this dist may vary for different pdf's $f = g$.

* **Example :**

$x_1, \dots, x_n \in \mathcal{X}$, $y_1, \dots, y_m \in \mathcal{Y}$ be two random samples of obj. represented in some data space \mathcal{X} . Suppose you have a f "

$d(x_i, y_j)$ that measures the distance between any two obj. $x_i, y_j \in X$. To test whether $x_1, \dots, x_n, y_1, \dots, y_m$ appear to come from the same distribution, the following test statistic is used

$$T_1 = \frac{2}{mn} \sum_{i=1}^n \sum_{j=1}^m d(x_i, y_j) - \frac{1}{\binom{n}{2}} \sum_{i \neq j} d(x_i, x_i) - \frac{1}{\binom{m}{2}} \sum_{i \neq j} d(y_i, y_i)$$

Hence, T_1 measures whether on average objects from the sample are more similar to each other than different samples.

Or we might consider a “nearest-neighbors” statistic: For each of the $m + n$ data values, look at the k other data values closest to it (as measured by the distance d) and count how many of these come from the same sample as itself. Let T_2 be the average of this count across all $m + n$ data points. So T_2 measures whether the k closest other objects tend to come from the same sample.

The distributions of T_1 and T_2 under H_0 may be difficult to understand theoretically and may depend on the unknown common distribution of $X_1, \dots, X_n, Y_1, \dots, Y_m$, but we can still carry out a permutation test based on T_1 or on T_2 .

A similar idea may be applied in the one-sample setting for testing the null hypothesis

$$H_0 : X_1, \dots, X_n \stackrel{\text{IID}}{\sim} f, \text{ for some PDF } f \text{ symmetric about 0}$$

based on the symmetry underlying the Wilcoxon signed-rank test.

* Generalised likelihood Ratio Test :

$$H_0 : \theta = \theta_0$$

$$H_1 : \theta \neq \theta_0$$

$$\Lambda = \frac{\text{lik}(\theta_0)}{\max_{\theta \in \Omega} \text{lik}(\theta)}$$

Ω is the space; our dist. is $\{(f(\theta)) | \theta \in \Omega\}$

Hypothesis Testing for categorical data : GLRT

1) Simple H_0 : let $\{\text{scenario}\}_{\theta \in \Theta}$ be a parametric model & $\theta_0 \in \Theta$ be a particular parameter value. We want to test:

$$H_0: \theta = \theta_0$$

$$H_1: \theta \neq \theta_0$$

The GLRT rejects for small values of the test statistic

$$\lambda = \frac{\text{lik}(\theta_0)}{\max_{\theta \in \Theta} \text{lik}(\theta)}$$

The numerator is the value of the likelihood at θ_0 where the denominator is the value of the likelihood at the MLE $\hat{\theta}$.

The level α -test rejects H_0 when $\lambda \leq c$ where c is chosen;

$$P_{H_0}(\lambda \leq c) \approx \alpha \quad (\text{probability of Type I error})$$

Note that the GLRT differs from the LRT in the context of the NP lemma, where the denominator was instead given by $\text{lik}(\theta_1)$ for the simple alternative $\theta = \theta_1$. Since the H_1 is not simple anymore the GLRT replaces the denominator by the maximum value of the likelihood over all values of θ .

* Problem: let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$

$$H_0: \theta = 0$$

$$H_1: \theta \neq 0$$

$$\prod_{i=1}^n \text{lik}(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2}\right) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n x_i^2\right)$$

$$\max_{\theta} \text{lik}(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x_i - \bar{x})^2}{2}\right) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2\right)$$

$$\begin{aligned} L &= \exp \left(\frac{1}{\alpha} \sum_{i=1}^n \left((x_i - \bar{x})^2 - \alpha x_i^2 \right) \right) \\ &= \exp \left(\frac{1}{\alpha} \sum_{i=1}^n (\bar{x}^2 - \alpha x_i \bar{x}) \right) \\ &= \exp \left(-\frac{n}{\alpha} \bar{x}^2 \right) \end{aligned}$$

$$\Rightarrow -2 \ln(L) = n \bar{x}^2 \quad (-L < c \text{ reject } H_0 \text{ or } -2 \ln(L) > c \text{ reject})$$

$$\Rightarrow -2 \ln(L) = n \bar{x}^2$$

Under H_0 : $\bar{x} \sim N(0, \frac{1}{n}) \Rightarrow n\bar{x} \sim N(0, 1)$
 $\Rightarrow n\bar{x}^2 \sim \chi^2_1$

$$\begin{aligned} -2 \ln(L) = n \bar{x}^2 &\stackrel{H_0}{\sim} \chi^2_1 \Rightarrow P_{H_0}(-2 \ln(L) > k) = \alpha \\ &\Rightarrow k = \chi^2_{1, 0.05} \text{ (upper } \alpha \text{ point)} \end{aligned}$$

* In general, the exact sampling distribution of $-2 \ln(L)$ under H_0 may not have a simple form as earlier but it may be approximated by a χ^2 dist. for large n .

* Theorem :

Let $f(x|\theta)$ be the parametric model & $x_1, \dots, x_n \stackrel{iid}{\sim} f(x|\theta_0)$. Suppose θ_0 is an interior point of Ω and the regularity conditions for consistency & asymptotic normality of the MLE hold. Then, $-2 \ln(L) \rightarrow \chi^2_k$ in distribution as $n \rightarrow +\infty$ where $k = \text{dimension of } \Omega$.

* Proof :

We will consider the case where $k = 1$ so θ is the single parameter. Let

$\ell(\theta)$ be the log-likelihood function and $\hat{\theta}$ be the MLE.

$$-L = \frac{\text{lik}(\theta_0)}{\text{lik}(\hat{\theta})}$$

Applying a Taylor series expansion of $\ell(\theta_0)$ around $\theta_0 = \hat{\theta}$

$$\begin{aligned}\ell(\theta_0) &\approx \ell(\hat{\theta}) + (\theta_0 - \hat{\theta}) \ell'(\hat{\theta}) + (\theta_0 - \hat{\theta})^2 \frac{\ell''(\hat{\theta})}{2} \\ &\approx \ell(\hat{\theta}) - \frac{1}{2} n I(\theta_0) (\theta_0 - \hat{\theta})^2\end{aligned}$$

$$\begin{aligned}\Rightarrow -2 \text{log}(-L) &= -2 \ell(\theta_0) + 2 \ell(\hat{\theta}) \\ &= n I(\theta_0) (\theta_0 - \hat{\theta})^2\end{aligned}$$

$\sqrt{n} I(\theta_0) (\hat{\theta} - \theta_0) \rightarrow N(0, 1)$ in distribution by asymptotic normality of MLE

$\Rightarrow \sqrt{n} I(\theta_0) (\hat{\theta} - \theta_0)^2 = -2 \text{log}(-L) \rightarrow \chi^2$, in distribution as $n \rightarrow +\infty$.

This theorem implies that an approximate level α test is given by rejecting H_0 when $-2 \text{log}(-L) > \chi_{\alpha}^2$, the upper α point of χ^2 dist.

* The dim K of \mathcal{R} is the # of free parameters in the model - the # of constraints. For instance, in the previous example there was only one parameter so the $\text{dim}(\mathcal{R}) = 1$. However, for a multinomial model p_1, \dots, p_k the $\text{dim} = K-1$ as there is a constraint $\sum_{i=1}^k p_i = 1$.

2) Ho composite : GLRT for testing a sub-model

More generally, let $\mathcal{S}_0 \subset \mathcal{R}$ be a subset of the parameter space \mathcal{R} corresponding to a lower dimensional submodel for testing

$H_0 : \theta \in \mathcal{S}_0$

$H_1 : \theta \notin \mathcal{S}_0 \quad \theta \in \mathcal{R}$

The GLRT statistic is defined as $\Lambda = \frac{\max_{\theta \in \Theta_0} \text{lik}(\theta)}{\max_{\theta \in \Theta} \text{lik}(\theta)}$.

In other words, Λ is the ratio of the likelihood evaluated at the MLE is the submodel and at the MLE in the full model.

For large n , under any $\theta_0 \in \Theta_0$, $-2\log(\Lambda)$ is approximately χ^2_k where k is the dimensionality of θ_0 and χ^2_k is an approximate level α -test rejects when $-2\log(\Lambda) > \chi^2_{k,\alpha}$ ($k = \dim(\theta_0) - \dim(\theta)$).

Example 6.2.1. (Hardy-Weinberg equilibrium⁹). At a single diallelic locus in the genome with two possible alleles A and a, any individual can have genotype AA, Aa, or aa. If we randomly select n individuals from a population, we may model the numbers of individuals with these genotypes as $(N_{AA}, N_{Aa}, N_{aa}) \sim \text{Multinomial}(n, (p_{AA}, p_{Aa}, p_{aa}))$.

When the alleles A and a are present in the population with proportions θ and $1 - \theta$, then under an assumption of random mating, quantitative genetics theory predicts that p_{AA} , p_{Aa} , and p_{aa} should be given by $p_{AA} = \theta^2$, $p_{Aa} = 2\theta(1 - \theta)$, and $p_{aa} = (1 - \theta)^2$ -this is called the Hardy-Weinberg equilibrium. In practice we do not know θ , but we may still test the null hypothesis that Hardy-Weinberg equilibrium holds for some θ :

$$H_0 : p_{AA} = \theta^2, p_{Aa} = 2\theta(1 - \theta), p_{aa} = (1 - \theta)^2 \text{ for some } \theta \in (0, 1).$$

This null hypothesis corresponds to a 1-dimensional sub-model (with a single free parameter θ) inside the 2-dimensional multinomial model (specified by general parameters p_{AA}, p_{Aa}, p_{aa} summing to 1). We may test H_0 using the GLRT:

The multinomial likelihood is given by

$$(N_{AA}, N_{Aa}, N_{aa}) \sim \text{Multinomial}(n, (p_{AA}, p_{Aa}, p_{aa}))$$

$$l(p_{AA}, p_{Aa}, p_{aa}) = \binom{n}{N_{AA}, N_{Aa}, N_{aa}} p_{AA}^{N_{AA}} p_{Aa}^{N_{Aa}} p_{aa}^{N_{aa}}. \quad \text{likelihood of } H_0$$

Letting $\hat{p}_{AA}, \hat{p}_{Aa}, \hat{p}_{aa}$ denote the full-model MLEs and $\hat{p}_{0,AA}, \hat{p}_{0,Aa}, \hat{p}_{0,aa}$ denote the sub-model MLEs, the generalized likelihood ratio is

$$\Lambda = \left(\frac{\hat{p}_{0,AA}}{\hat{p}_{AA}} \right)^{N_{AA}} \left(\frac{\hat{p}_{0,Aa}}{\hat{p}_{Aa}} \right)^{N_{Aa}} \left(\frac{\hat{p}_{0,aa}}{\hat{p}_{aa}} \right)^{N_{aa}},$$

so In order to apply the theorem

$$-2 \log \Lambda = 2N_{AA} \log \frac{\hat{p}_{AA}}{\hat{p}_{0,AA}} + 2N_{Aa} \log \frac{\hat{p}_{Aa}}{\hat{p}_{0,Aa}} + 2N_{aa} \log \frac{\hat{p}_{aa}}{\hat{p}_{0,aa}}. \quad (1)$$

↑ ↗ # of individuals with this gene

The full-model MLEs are given by $\hat{p}_{AA} = N_{AA}/n$, $\hat{p}_{Aa} = N_{Aa}/n$, and $\hat{p}_{aa} = N_{aa}/n$, by

Example 1.3.4 from Chapter 2. To find the sub-model MLEs, note that under H_0 , the multinomial likelihood as a function of θ is

$$\text{lik}(\theta) = \binom{n}{N_{AA}, N_{Aa}, N_{aa}} (\theta^2)^{N_{AA}} (2\theta(1-\theta))^{N_{Aa}} ((1-\theta)^2)^{N_{aa}}$$

$$= \binom{n}{N_{AA}, N_{Aa}, N_{aa}} 2^{N_{Aa}} \theta^{2N_{AA}+N_{Aa}} (1-\theta)^{N_{Aa}+2N_{aa}}.$$

Maximizing the likelihood over parameters (p_{AA}, p_{Aa}, p_{aa}) belonging to the sub-model is equivalent to maximizing the above over θ . Differentiating the logarithm of the above likelihood and setting it equal to 0, we obtain the MLE

$$\hat{\theta} = \frac{2N_{AA} + N_{Aa}}{2N_{AA} + 2N_{Aa} + 2N_{aa}} = \frac{2N_{AA} + N_{Aa}}{2n}$$

for θ , which yields the sub-model MLEs

$$\hat{p}_{0,AA} = \left(\frac{2N_{AA} + N_{Aa}}{2n} \right)^2$$

$$\hat{p}_{0,Aa} = 2 \left(\frac{2N_{AA} + N_{Aa}}{2n} \right) \left(\frac{N_{Aa} + 2N_{aa}}{2n} \right)$$

$$\hat{p}_{0,aa} = \left(\frac{N_{Aa} + 2N_{aa}}{2n} \right)^2.$$

Substituting these expressions into equation (1) yields the formula for $-2 \log \Lambda$ in terms of the observed counts N_{AA}, N_{Aa}, N_{aa} . The difference in dimensionality of the two models is $2 - 1 = 1$, so an approximate level- α test would reject H_0 when $-2 \log \Lambda$ exceeds $\chi_1^2(\alpha)$.
 2 : full model - $\Delta \log \Lambda > \chi_1^2(\alpha)$ reject.
 1 : submodel

* Test of Independence :

		Dem	Rep	Independent	
Gender	Male	422	273	381	$P_{1,0}$
	Female	299	232	365	$P_{2,0}$
		$P_{1,1}$	$P_{2,1}$	$=$	"
					$N = 1972$

In this example the following table cross-classifies a random sample of 1972 people by gender and by political party identification :

39% of females identified as democrat and 25% as republican while 33% of males identified as democrat and 26% as republican. Is there an evidence of an association b/w gender & party identification from which the sample was drawn.

Let N_{ij} be the number of people with gender i who voted for party j : $\sum_{j=1}^3 \sum_{i=1}^2 N_{ij} = n$ is the # of observations

Let P_{ij} be the probability of belonging to (i,j) th group. We may count N_{ij} as multinomial distribution with prob p_{ij} ($\sum_i \sum_j p_{ij} = 1$)

We can write the total no of people in group i : $\sum_{j=1}^3 N_{ij} = N_{i,0}$

The prob : $P_{i,0} = \sum_j p_{ij}$ prob of belonging to gender i
 $P_{0,j} = \sum_i p_{ij}$ prob of belonging to party j

Test : $H_0 : P_{ij} = P_{i,0} P_{0,j} \quad \forall i, j$ Reject \Rightarrow there is an association

Under H_0 , $\dim(\text{submodel}) = 5 - 2 = 3$ and the $\dim(\text{full model}) = 6 - 1 = 5$. For the GLRT has an approximately null distribution with $5 - 3 = 2$ degrees of freedom.

To derive the form :

$$\begin{aligned} \mathcal{L}_{\text{sub}}(x) &= \log \left(\binom{n}{N_1, \dots, N_K} \prod_{i=1}^K p_i^{N_i} \right) \quad (N_1, \dots, N_K) \sim (n(p_1, \dots, p_K)) \\ &= \log \left(\binom{n}{N_1, \dots, N_K} \right) + \sum_{i=1}^K N_i \log(p_i) \end{aligned}$$

Letting $\hat{P}_{0,i}$ denote the MLEs in the submodel π_0 and \hat{p}_i denote the MLEs in the full model.

$$\text{GLRT} : \Lambda = \prod_{i=1}^K \left(\frac{P_{0,i}}{\hat{p}_i} \right)^{N_i} \Rightarrow -2 \log(\Lambda) = -2 \sum_{i=1}^K N_i \log \left(\frac{P_{0,i}}{\hat{p}_i} \right)$$

Since $\hat{P}_i = \frac{N_i}{n}$, let $E_i = \hat{P}_{0..i} \cdot n$ denote the expected count for outcome i corresponding to Ω : $-2\text{log}(n) = -2 \sum_{i=1}^k N_i \text{log}\left(\frac{N_i}{E_i}\right)$

Under H_0 : $P_{i,j} = P_{i,0} P_{0,j} \quad \forall i, j$

$$\begin{aligned} \text{lik}(P_{1,0}, P_{2,0}, P_{0,1}, P_{0,2}, P_{0,3}) &= \binom{n}{N_{1,0}, \dots, N_{2,3}} \prod_{i=1}^2 \prod_{j=1}^3 (P_{i,0} P_{0,j})^{N_{i,j}} \\ &= \binom{n}{N_{1,0}, \dots, N_{2,3}} \prod_{i=1}^2 P_{i,0}^{N_{i,0}} \prod_{j=1}^3 P_{0,j}^{N_{0,j}} \end{aligned}$$

$$\begin{aligned} \ell(P_{1,0}, \dots, P_{0,3}) &= \text{log} \left(\binom{n}{N_{1,0}, \dots, N_{2,3}} \right) + \sum_{i=1}^2 N_{i,0} \text{log}(P_{i,0}) + \sum_{j=1}^3 N_{0,j} \text{log}(P_{0,j}) \\ &\quad + A \left(\sum_{i=1}^2 P_{i,0} - 1 \right) + B \left(\sum_{j=1}^3 P_{0,j} - 1 \right) \end{aligned}$$

$$\begin{aligned} \text{Select } \lambda = -n \text{ and } \gamma = -n \Rightarrow P_{i,0} &= \frac{-N_{i,0}}{\lambda} = \frac{N_{i,0}}{n} \\ P_{0,j} &= \frac{-N_{0,j}}{\gamma} = \frac{N_{0,j}}{n} \end{aligned}$$

$$\begin{aligned} \hat{P}_{i,j} &= \left(\frac{N_{i,0}}{n} \right) \left(\frac{N_{0,j}}{n} \right) \Rightarrow E_{i,j} = n \hat{P}_{i,j} \\ &\Rightarrow -2\text{log} n = 8.31 \\ &\text{p-value} = 0.014 \text{ Reject!} \end{aligned}$$

Hence, there is an association between gender & political identification

* Test of Homogeneity:

We have two independent count obs from 2 multinomial samples and each have K outcomes $(N_1, \dots, N_K) \sim MN(n, (p_1, \dots, p_K))$ and $(M_1, \dots, M_K) \sim MN(m, (q_1, \dots, q_K))$ where n and m are known sample sizes. To test the homogeneity $H_0: P_i = q_i \quad \forall i \in \{1, \dots, K\}$

The following table counts the occurrences of six different short words in Chapter 1 and Chapter 6 of Semelvren written by Austin and in Chapters 12 and 24 of Semelction written by the examiner.

	a	an	ohns	that	with	without
Ch 1 & Ch 6	101	11	15	37	28	10 202
Ch 12 & Ch 24	83	29	15	22	43	4 194

Let us model the counts from Ch 1 & Ch 6 $\sim \text{MV}(\mu_1, (\rho_1, \dots, \rho_6))$ and Ch 12 & Ch 24 $\sim \text{MV}(\mu_2, (\rho_1, \dots, \rho_6))$. We wish to test $H_0: \rho_i = q_i \forall i \in \{1, \dots, 6\}$.

To derive the GLRT:

$$\text{lik}(\rho_1, \dots, \rho_K, q_1, \dots, q_K) = \binom{n}{N_1, \dots, N_K} \prod_{i=1}^K \rho_i^{N_i} \binom{m}{M_1, \dots, M_K} \prod_{i=1}^K q_i^{M_i}$$

Let $\hat{\rho}_i = \frac{N_i}{n}$, $\hat{q}_i = \frac{M_i}{m}$ be the full model MLEs & $\hat{\rho}_{0,i} = \hat{q}_{0,i}$ be the MLEs of the submodel.

$$\Lambda = \prod_{i=1}^K \left(\frac{\rho_{0,i}}{\hat{\rho}_i} \right)^{N_i} \prod_{i=1}^K \left(\frac{q_{0,i}}{\hat{q}_i} \right)^{M_i} \Rightarrow$$

$$-\lambda \text{log}(\Lambda) = -\sum_{i=1}^K (N_i \text{log} \left(\frac{\hat{\rho}_i}{\rho_{0,i}} \right) + M_i \text{log} \left(\frac{\hat{q}_i}{q_{0,i}} \right))$$

$$= -\sum_{i=1}^K (N_i \text{log} \left(\frac{\hat{\rho}_i}{E_i} \right) + M_i \text{log} \left(\frac{\hat{q}_i}{F_i} \right))$$

$$E_i = n \hat{\rho}_i \quad F_i = m \hat{q}_i$$

$$\text{In the submodel : } \text{lik}(\rho_1, \dots, \rho_K) = \binom{n}{N_1, \dots, N_K} \binom{m}{M_1, \dots, M_K} \prod_{i=1}^K \rho_i^{N_i+M_i}$$

$$\begin{aligned} \text{log lik}(\rho_1, \dots, \rho_K) &= \text{log}(\alpha) + \text{log}(\alpha') + \sum_{i=1}^K (N_i+M_i) \text{log} \rho_i \\ &\quad + \lambda \left(\sum_{i=1}^K \rho_i - 1 \right) \end{aligned}$$

$$\Rightarrow \rho_i = \frac{M_i+N_i}{-\lambda} \quad \text{choose } \lambda = -(n+m)$$